# MLM Training Example:

Original: "The cat sits on the mat"



| | The | [MASK] | sits | on | the | [MASK] |
|---|---|---|---|---|---|---|
| The | | | | | | |
| [MASK] | | | | | | |

pos 0    pos 1    pos 2    pos 3    pos 4    pos 5

0.25   0.35   0.15   0.20   0.05

Predicts: "cat"      Predicts: "mat"
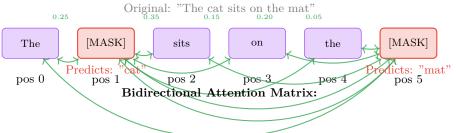
**Bidirectional Attention Matrix:**

**Key Differences from Autoregressive Models**

**Bidirectional Context:**    **No Causal Masking:**    **Parallel Prediction:**

Uses both left & right context   Future tokens inform predictions   Multiple masks predicted jointly

**MLM Objective:** Predict original tokens using full bidirectional context